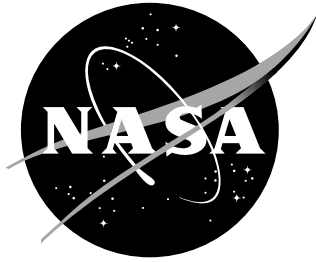


NASA/TM-1998-208956



# Architectural Optimization of Digital Libraries

*Aileen O. Biser*  
*Langley Research Center, Hampton, Virginia*

---

December 1998

## The NASA STI Program Office ... in Profile

Since its founding, NASA has been dedicated to the advancement of aeronautics and space science. The NASA Scientific and Technical Information (STI) Program Office plays a key part in helping NASA maintain this important role.

The NASA STI Program Office is operated by Langley Research Center, the lead center for NASA's scientific and technical information. The NASA STI Program Office provides access to the NASA STI Database, the largest collection of aeronautical and space science STI in the world. The Program Office is also NASA's institutional mechanism for disseminating the results of its research and development activities. These results are published by NASA in the NASA STI Report Series, which includes the following report types:

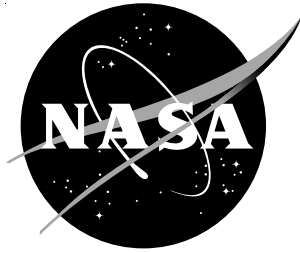
- **TECHNICAL PUBLICATION.** Reports of completed research or a major significant phase of research that present the results of NASA programs and include extensive data or theoretical analysis. Includes compilations of significant scientific and technical data and information deemed to be of continuing reference value. NASA counterpart of peer-reviewed formal professional papers, but having less stringent limitations on manuscript length and extent of graphic presentations.
- **TECHNICAL MEMORANDUM.** Scientific and technical findings that are preliminary or of specialized interest, e.g., quick release reports, working papers, and bibliographies that contain minimal annotation. Does not contain extensive analysis.
- **CONTRACTOR REPORT.** Scientific and technical findings by NASA-sponsored contractors and grantees.
- **CONFERENCE PUBLICATION.** Collected papers from scientific and technical conferences, symposia, seminars, or other meetings sponsored or co-sponsored by NASA.
- **SPECIAL PUBLICATION.** Scientific, technical, or historical information from NASA programs, projects, and missions, often concerned with subjects having substantial public interest.
- **TECHNICAL TRANSLATION.** English-language translations of foreign scientific and technical material pertinent to NASA's mission.

Specialized services that complement the STI Program Office's diverse offerings include creating custom thesauri, building customized databases, organizing and publishing research results ... even providing videos.

For more information about the NASA STI Program Office, see the following:

- Access the NASA STI Program Home Page at <http://www.sti.nasa.gov>
- E-mail your question via the Internet to [help@sti.nasa.gov](mailto:help@sti.nasa.gov)
- Fax your question to the NASA STI Help Desk at (301) 621-0134
- Phone the NASA STI Help Desk at (301) 621-0390
- Write to:  
NASA STI Help Desk  
NASA Center for AeroSpace Information  
7121 Standard Drive  
Hanover, MD 21076-1320

NASA/TM-1998-208956



# Architectural Optimization of Digital Libraries

*Aileen O. Biser*  
*Langley Research Center, Hampton, Virginia*

National Aeronautics and  
Space Administration

Langley Research Center  
Hampton, Virginia 23681-2199

---

December 1998

---

Available from:

NASA Center for AeroSpace Information (CASI)  
7121 Standard Drive  
Hanover, MD 21076-1320  
(301) 621-0390

National Technical Information Service (NTIS)  
5285 Port Royal Road  
Springfield, VA 22161-2171  
(703) 605-6000

## **ABSTRACT**

### **ARCHITECTURAL OPTIMIZATION OF DIGITAL LIBRARIES**

Aileen O. Biser

Old Dominion University, 1998

Co-Chairs of Advisory Committee: Dr. Kurt Maly

Dr. Stewart N. T. Shen

This work investigates performance and scaling issues relevant to large scale distributed digital libraries. Presently, performance and scaling studies focus on specific implementations of production or prototype digital libraries. Although useful information is gained to aid these designers and other researchers with insights to performance and scaling issues, the broader issues relevant to very large scale distributed libraries are not addressed. Specifically, no current studies look at the extreme or worst case possibilities in digital library implementations. A survey of digital library research issues is presented. Scaling and performance issues are mentioned frequently in the digital library literature but are generally not the focus of much of the current research.

In this thesis a model for a Generic Distributed Digital Library (GDDL) and nine cases of typical user activities are defined. This model is used to facilitate some basic analysis of scaling issues. Specifically, the calculation of Internet traffic generated for different configurations of the study parameters and an estimate of the future bandwidth needed for a large scale distributed digital library implementation.

This analysis demonstrates the potential impact a future distributed digital library implementation would have on the Internet traffic load and raises questions concerning

the architecture decisions being made for future distributed digital library designs and the Internet capacities that will be necessary to support them. This analysis suggests that network capacities of 622 Mbps will be required to go much beyond 100 heavily used independent digital library sites. Additionally, capacities beyond 622 Mbps will be required to realize the worldwide distributed digital library consisting of a 1000 or more digital library sites. These results also point out the need for architecture modifications and software improvements to reduce and minimize the amount of network traffic generated as we move to a global digital library implementation.

## ACKNOWLEDGMENTS

Professors Kurt Maly and Stewart N. T. Shen provided the direct advisement for this research.

NASA Langley Research Center has provided me with the opportunity and resources to perform digital library research. In particular, I would like to thank Michael Little and Mary McCaskill for allowing me the time and solitude to focus on this endeavor. I would like to thank Michael Nelson for the support, encouragement and guidance during the past two years that has made this possible. I would like to thank Frank Thames for encouraging me to continue, and Cathy Cronin for listening and keeping my spirits positive.

Finally, I would like to thank my husband and biggest supporter, Jerry, who always believes I can accomplish anything even when I doubt myself, and my sons, Aaron and Evan, who bring joy to my life and make every day meaningful. The work we do today will change the world for our children in many wonderful ways.





## TABLE OF CONTENTS

	PAGE
LIST OF TABLES.....	x
LIST OF FIGURES.....	xii
 Section	
1. INTRODUCTION.....	1
2. BRIEF REVIEW OF DIGITAL LIBRARIES.....	3
2.1 Digital library definition.....	3
2.2 The future of digital libraries.....	3
2.3 Definition of a distributed digital library.....	4
2.4 Examples of current distributed digital libraries.....	6
3. STATE OF ART IN DIGITAL LIBRARY RESEARCH.....	8
3.1 Survey of digital library research.....	8
3.2 Survey of digital library performance research.....	9
3.3 Discussion of digital library simulation studies.....	10
4. PROBLEM STATEMENT.....	15
4.1 Statement of the study question.....	15
4.2 Justification.....	16
4.3 Analysis and comparison of simulation studies.....	17
4.4 Discussion.....	19

	PAGE
5. PROBLEM ANALYSIS.....	21
5.1 Generic model design.....	21
5.1.1 Description of system components.....	21
5.1.2 Description of model data flow.....	24
5.2 Model specifications.....	25
5.3 Examples of digital libraries.....	28
5.3.1 Physics E-Print.....	28
5.3.2 NTRS.....	30
5.3.3 NCSTRL.....	32
5.4 Parameters.....	35
5.5 Measurements and supporting data.....	35
5.6 Discussion of study cases.....	38
5.6.1 Breakdown of cases studied.....	41
5.6.2 Case analysis.....	42
5.7 Study assumptions.....	44
5.8 Analytic formulas.....	46
5.9 Results tables.....	48
6. FINDINGS AND INTERPRETATIONS.....	52
7. FUTURE RESEARCH.....	56
8. SUMMARY AND CONCLUSIONS.....	58

	PAGE
8.1 Summary of contributions.....	58
8.2 Conclusions.....	59
REFERENCES.....	60

## LIST OF TABLES

TABLE	PAGE
1. Distribution of digital libraries.....	4
2. Current distributed digital libraries.....	7
3. Digital library performance and scaling studies.....	10
4. Primary goal of the studies.....	17
5. Model components defined.....	18
6. Measurements used in the studies.....	18
7. Differences in study implementations.....	18
8. Parameters varied in the studies.....	19
9. Model nomenclature.....	21
10. GDDL model component specifications.....	27
11. Example digital libraries.....	28
12. NCSTRL specific components.....	33
13. Primary model parameters.....	35
14. Internet technology.....	36
15. Internet throughputs.....	36
16. Average values measured from LTRS.....	37
17. User session characteristics.....	39
18. Digital library usage cases.....	40
19. Case breakdown by percentages.....	41

TABLE	PAGE
20. Case breakdown by user count.....	42
21. Equations used for case analysis.....	44
22. Total traffic generated per individual case for worst case analysis.....	49
23. Total traffic generated per individual case for average analysis.....	49
24. Calculation of total traffic for Worst Case using Sample A .....	50
25. Calculation of total traffic for Average Case using Sample A.....	50
26. Calculation of total traffic for Worst Case using Sample B.....	51
27. Calculation of total traffic for Average Case using Sample B.....	51
28. Time to transmit at 130 Mbps.....	53
29. Time to transmit at 450 Mbps.....	54

## LIST OF FIGURES

FIGURE	PAGE
1. Generic Distributed Digital Library model.....	22
2. Local data flow of GDDL.....	24
3. Global data flow of GDDL.....	25
4. Physics E-Print model.....	29
5. Global data flow of Physics E-Print.....	30
6. NTRS model.....	31
7. Local data flow of NTRS.....	31
8. Global data flow of NTRS.....	32
9. NCSTRL model.....	33
10. Local data flow of NCSTRL.....	34
11. Global data flow of NCSTRL.....	34

## SECTION ONE

### INTRODUCTION

The field of digital library research is young, broad and growing rapidly. The problems yet to be solved cross the entire spectrum of computer science, information science, human-computer interaction, publishing and commercialization. Research is simultaneously occurring in many different areas all with the effort to develop or improve a digital library for many users. What happens when these digital library efforts and many others come to pass and we have access to hundreds of digital libraries? This is the primary focus of this study. Specifically, we would like to determine the Internet traffic that can be anticipated in the future with hundreds and possibly thousands of digital libraries available to the world users.

The approach to solving this problem is to define the basic components of a distributed digital library (DDL) and use that knowledge to perform further high level analysis of a DDL independent of any specific implementation issues. It is suggested that by using this basic set of components the function of a DDL can be represented, analyzed, and simulated in order to obtain insight into architecture changes beneficial in a broad sense. By defining the basic components and suggesting a typical user usage pattern, we have the basic elements necessary to express architecture and usage pattern changes. This will allow for the calculation and analysis of these changes. The results

obtained will show that for at least the lower bound worst case analysis Internet traffic will indeed be a large problem for growth beyond 100 heavily used distributed digital library sites on the Internet.

The outline for the rest of this thesis is as follows: Section two provides a brief review of digital libraries with a definition and examples of distributed digital libraries. Section three provides a survey of digital library research with examples of distributed digital libraries and a look at the performance and simulation studies that have been done. Section four formally defines the problem to be solved and provides a justification for the work. Section five presents the main analysis and discusses the Generic Distributed Digital Library model and nomenclature, presents representations of other digital libraries using this nomenclature, defines cases of user activities that will be used in the total traffic calculation and finally presents the formulas and results obtained. Section six discusses the Internet traffic calculations and impact of these findings. Section seven discusses the limitations of this study and the future work needed to improve the validity and accuracy of the results. We conclude with Section eight.



## **SECTION TWO**

### **BRIEF REVIEW OF DIGITAL LIBRARIES**

#### **2.1 Digital library definition**

The term digital library causes much confusion in general conversation. Depending on an individual background and the context in which the term is used, each person may assume something different. For purposes of this thesis we will define a digital library according to "Digital Libraries are organized collections of digital information" (Lesk 1997).

#### **2.2 The future of digital libraries**

As Lesk also points out, individuals or groups that select, organize and catalog large numbers of pages have turned the World Wide Web into many Digital Libraries. It is obvious from a survey of the literature that many and diverse digital libraries are being developed. The future will be populated with many digital libraries but what that future really looks like is partly speculation and assumptions based on current examples. What we do know is that digital libraries are here to stay in possibly many forms and hopefully will be integrated for ease of use.

One specific example of a future digital library is NCSTRL+ (Nelson et al. 1998). This is an important example of the direction some digital library research is taking by providing access to information and its associated parts, be they data, software, graphics or video. It is fair to say that the digital library of the future will provide not only access

to documents, but to all types of data in some logical and user friendly fashion. This is important to note because this study is limited in its ability to analyze future digital library architecture issues because the data needed does not exist. Data available today and used in this analysis is only representative of the current limited implementations of digital libraries. As a result, many assumptions and projections of possibilities are made.

### **2.3 Definition of a distributed digital library**

Taxonomies in Digital Libraries have been studied (Esler and Nelson 1998) and this early work resulted in the definition of a nomenclature for describing various digital library projects. They can be differentiated by their architecture (distributed or centralized) and by the identity of the sponsor of the digital library (traditional publishers or authoring individuals/groups). These four major architectural categories for identifying Digital Libraries established by Esler and Nelson are shown in Table1.

**Table 1.** Distribution of digital libraries

	Distributed	Centralized
Traditional Publisher	DP	CP
Authoring Individual/ Organization	DO	CO

Esler and Nelson give us the following definitions:

*“Centralized Architecture, Traditional Publisher (CP) - Input is from traditional publishing sources such as journals and professional societies, and all input is collected in a single physical and logical location. The server is either up or down, there is no graduated level of availability....”*

*“Distributed Architecture, Traditional Publisher (DP) - Input is from traditional publishing sources such as journals and professional societies, but the input is not transmitted to a single physical location. The user interface may give the appearance of a central location, but the service is comprised of several servers....”*

*“Centralized Architecture, Authoring Individual/Organization (CO) - Input is from either individuals (a few papers at a time) or from an organization (papers transmitted in batches) and the input is transferred to a central location for indexing, processing and redistribution....”*

*“Distributed Architecture, Authoring Individual/Organization (DO) - Input could still be from individuals, but separate servers encourage clustering of publishers along organizational boundaries. Input stays at the server to which it was posted and the user interface handles querying all appropriate servers and collating and presenting the results....”*

From a performance and scaling perspective where we are looking at issues of network traffic and communication load, these four classifications can be more narrowly defined as either distributed or centralized. A distributed digital library is characterized as having multiple services distributed throughout an Internet and/or Intranet. In this

architecture the user has access either locally via an Intranet to a subset of the digital library services or access globally via the Internet to all or a broadly defined subset of the digital library services. In a centralized digital library a single point of access provides services to a local or distributed user community. In the centralized case the network traffic is characterized by many users from many locations (Internet or Intranet) accessing a single server providing all digital library services. This is contrasted with the network characteristics of a distributed digital library where many users communicate with many distinct services distributed globally and locally. In terms of network traffic measurements and analysis, the distributed digital library is many times more complex to analyze than in the case of a centralized digital library.

As pointed out in (Esler and Nelson 1998), these classification factors are important because it is suggested that distributed architecture digital libraries are more likely to be scalable than centralized digital libraries.

## **2.4 Examples of current distributed digital libraries**

Table 2 provides examples of current production and prototype distributed digital libraries. The limitation that was placed on inclusion in this example set is that the digital library architecture conforms to our definition of a distributed digital library stated in Section 2.3. In surveying available digital libraries we find that many WWW accessible digital libraries (Nelson 1998) have centralized archives and are therefore not represented in Table 2.

**Table 2.** Current distributed digital libraries

<b>DL Identifier</b>	<b>DL Name and URL</b>	<b>Content</b>
DLI	Digital Library Initiative Not available to the public <a href="http://dli.grainger.uiuc.edu">http://dli.grainger.uiuc.edu</a>	Multi-discipline
NTRS	NASA Technical Report Server <a href="http://techreports.larc.nasa.gov/cgi-bin/ntrs">http://techreports.larc.nasa.gov/cgi-bin/ntrs</a>	NASA technical reports
NCSTRL	Network Computer Science Technical Report Library <a href="http://www.ncstrl.org">http://www.ncstrl.org</a>	Computer science technical reports
NCSTRL+	Experimental and in development <a href="http://dlib.cs.odu.edu">http://dlib.cs.odu.edu</a>	Multi-discipline, multi-format data objects
UCSTRI	Unified Computer Science Technical Report Index <a href="http://www.cs.indiana.edu/cstr/search">http://www.cs.indiana.edu/cstr/search</a> (VanHeyningen 1994)	Computer science technical reports
NIX	NASA Image Exchange <a href="http://nix.nasa.gov">http://nix.nasa.gov</a> (von Ofenheim et al. 1998)	NASA videos and images
EOSDIS	Earth Observing System Data and Information System <a href="http://www-v0ims.gsfc.nasa.gov/v0ims/eosdis_home.html">http://www-v0ims.gsfc.nasa.gov/v0ims/eosdis_home.html</a>	Satellite data and related products
ADS	Astrophysics Data System <a href="http://ads.harvard.edu">http://ads.harvard.edu</a> (Eichhorn 1998)	Astrophysics and related technical documents
Arquitect	Portuguese National Digital Library (Borbinha et al. 1997)	Multi-document classifications
Medoc	German digital library project <a href="http://medoc.informatik.uni-hamburg.de">http://medoc.informatik.uni-hamburg.de</a> (Adler et al. 1998)	Technical reports, grey literature and multi collections
NHSE	National HPCC Software Exchange <a href="http://www.nhse.org">http://www.nhse.org</a> (Browne et al. 1995)	High performance and parallel computing software, documents, data and information

## **SECTION THREE**

### **STATE OF ART IN DIGITAL LIBRARY RESEARCH**

#### **3.1 Survey of digital library research**

A survey of the current digital library research shows that much of the effort is focused on creating testbed digital libraries with emphasis on infrastructure (Lynch and Garcia-Molina 1995; Nurnberg et al. 1995; Chen et al. 1996), protocols (Gravano et al. 1997a), indexing (Esler and Nelson 1997), federation (Shatz et al. 1996), digital objects (Kahn and Wilenski 1995; Lagoze and Ely 1995), and interoperability (Maa et al. 1997). Today's primary research goal is to build the digital library of the future with attempts to create large enough testbeds to do further research on the issues of scaling. It is widely agreed that scaling is a critical research issue in developing large-scale digital libraries (Shatz and Chen 1996). However, this is considered a deep research problem, which requires the deployment of large-scale systems for experimentation. At this time there exist substantial functional digital libraries (such as NTRS and NCSTRL) that are used daily and growing. These existing systems have already faced performance and design issues (Nelson and Maa 1996; French 1996; Balci et al. 1998; French et al. 1998) as they grow and evolve. It is clear that performance scaling analysis and tuning of architectural choices are issues that should be addressed today. The examination of functioning digital library projects and current research efforts reveals that there are a number of distinct architectural approaches to building digital libraries (Esler and Nelson 1998). A closer

examination and analysis of these approaches should provide insight into which approaches are expected to scale well as we move toward large-scale digital library systems.

We suggest that the problems of scaling and performance must be evaluated today for systems in use and new design options being considered. In evaluating these problems we will lay the groundwork for optimization of future digital library architectures.

### **3.2 Survey of digital library performance research**

In researching the issues of performance and scaling in digital libraries a number of different studies were identified that in some way addressed these issues and are shown in Table 3. The primary focus of the studies varied greatly from query optimization to server utilization issues and the approach used to address the questions was also varied. Of the various studies conducted only two incorporated a simulation of the system to experiment with and analyze architecture changes. We discuss these two studies in detail in the next Section.

**Table 3.** Digital library performance and scaling studies

<b>DL Name</b>	<b>Reference</b>	<b>Approach</b>	<b>Primary Focus</b>
NTRS	Nelson and Maa 1996	Data analysis and software modification	Parallel searches to reduce query response time
NTRS	Esler and Nelson 1997	Testing and data analysis	Development of NASA indexing benchmarks and results
NCSTRL	French 1996	Model analysis	Query processing time and performance bottlenecks
NCSTRL	French et al. 1998	Data analysis	Query routing to reduce distributed search time
NCSTRL	Balci et al. 1998b	Simulation	General performance analysis tool
INQUERY	Cahoon and Mckinley 1995; 1996; 1997	Prototype system and simulation analysis	Analyze effect of scaling to multiple servers
DLI	McGrath 1996	Interviews and analysis	Evaluation of scaling issues
ADL	Andresen et al. 1996	Prototype system analysis	Network bandwidth requirements and computational and I/O demands
STARTS	Gravano et al. 1997b	Data Analysis	Performance of payment schemes
KEYNET	Baclawski 1995	Prototype system analysis	Scalability of distributed information retrieval queries

### 3.3 Discussion of digital library simulation studies

The paper (Balci et al. 1998b) describes the design of a simulation of NCSTRL using the VSE (Visual Simulation Environment) (Balci et al. 1998a). A number of reusable model components were defined for NCSTRL to be configurable in the simulation. These



components were defined with the capabilities of the Dienst 4.0 architecture of the NCSTRL implementation.

The components defined include Top Level, Region, Dienst Server (simulates distributed searches) (Lagoze et al. 1995), Merged Index Server, Central Index Server, Backup Server, User Population (models submission of queries to a particular server), and Query. The workload characterization simulated includes query integration time, server response to queries and transaction time of request. Log data from three servers was used to characterize these times.

This model simulates Dienst 4.0 (Davis and Lagoze 1994; Davis et al. 1995) version of NCSTRL (Davis and Lagoze 1996) and does not represent the current architecture, NCSTRL 5.0 and Dienst 4.1. In order to simulate NCSTRL as it is today the function of different model components would have to be modified and/or new model components defined. The simulation of users as a User Population is unclear and the paper does not fully describe this component. It appears to assume that all users interface first to their local Dienst server user interface and not to the main top-level user interface. User Population queries go first to the local Dienst server and from there to the Region server and beyond.

This paper does not present any results of the simulations and gives few details of the input parameters available to the users. It does state that the user can run the simulation interactively to observe the actions of the architecture changes being simulated or in background mode to collect statistical information for later analysis.

This study differs from ours in a number of ways. First it is a simulation of a specific digital library implementation (NCSTRL) and a specific architectural implementation of that digital library frozen in time. The results produced visually in an interactive fashion or statistically serve to assist decision-making concerning the NCSTRL architecture only. No suggestion is made that results from this simulation can be used to assist other digital library designers or implementers in making decisions concerning their architectural choices. Some general knowledge can be gained from the results but no clear guidance can be derived for other digital library implementations.

Additional architecture and simulation studies were done by Cahoon and McKinley at the University of Massachusetts. The basis of this work began with an analysis (Cahoon and McKinley 1995) of a prototype distributed information retrieval system based on Inquiry (Callan et al. 1992), an existing, unified Information Retrieval system. This study continued (Cahoon and McKinley 1996; Cahoon and McKinley 1997) with the development of a simulation to conduct workload analysis of the prototype distributed Inquiry system. These studies were conducted to determine if the Inquiry Information Retrieval Server could be distributed across multiple systems and maintain acceptable service. Acceptable service is determined by observed response time degradation and increased system utilization of the servers.

Although this study does not refer to this architecture as a digital library system, it is included here because we feel that the architecture and components conform to the definition of a digital library as defined in Section 2.1. The system consists of Inquiry

servers, a connection server and clients. The study focused on the development of a distributed prototype and a simulation of that prototype. Data used in the simulation for workload analysis and parameter values were obtained either from the operational distributed Inquiry prototype or a production Inquiry system. The workload characterization for this simulation included: Query Evaluation Time, Document Retrieval Time, Summary Retrieval Time, Connection Server Time, Time to Merge Results, and Network Time. The system parameters that are varied in the study include the number of users, size and total number of documents in the collections, terms per query, query term frequency, user think time, number of answers returned, and workload.

The study examined distributing a single Inquiry text collection across multiple systems and the management of multiple distinct text collections on independent servers. In both cases a single central broker (or connection server) was used to interface between the users and the individual Inquiry servers. Much of the emphasis was on varying information retrieval parameters such as terms per query, user think time and document collection sizes. Network time was limited to sender and receiver overhead and network latency on a 10Mbps Ethernet LAN. A number of tests were conducted varying the simulation parameters and the results evaluated based on average transaction sequence time, connection server utilization and Inquiry server utilization. For many of the configurations tested, the connection server was the bottleneck to performance. The study is useful in presenting the bottlenecks and usage patterns that lead to the best response time and system utilization for an Inquiry implementation. We can gain some

insight into how other implementations may act in similar configurations. For this study the connection server was identified as a limiting factor for scaling and suggestions were given to correct this problem. This is consistent with the study done by (Fuhr 1997) which points out the need for multiple brokers in networked information retrieval of multiple data sources.

This study differs from ours in one very important way. The Inquiry study only takes into consideration local area network traffic where our study is mainly interested in wide area network traffic. Our focus is on the impact multiple digital libraries have on wide area traffic, while the Inquiry study focused on the ability of the connection server and Inquiry servers to respond to different workloads and configurations.

## **SECTION FOUR**

### **PROBLEM STATEMENT**

Implementers of Digital Libraries today and in the future will be faced with architectural design decision that will be difficult to make without the help of performance and scaling data from production implementations, testbed research and simulation studies.

The objective of this project is to investigate the design and performance characteristics of digital library architectures and the scaling issues critical for the design of an optimum large-scale distributed digital library. This research can be facilitated by studying the architectural approaches that have been implemented in existing functional digital libraries such as the Physics E-Print Digital Library (Ginsbarg 1994), the NASA Technical Report Server (NTRS) (Nelson et al. 1995), or the Network Computer Science Technical Report Library (NCSTRL) (Davis and Lagoze 1996).

The approach used in this research is to conduct an analytic study of a generic distributed digital library architecture with emphasis on the performance and scaling issues relevant to future digital libraries. The results of this study will facilitate the ongoing research to design large-scale digital library architectures and assist in making design decisions for existing functional digital libraries.

#### **4.1 Statement of the study question**

The main focus of this study is to determine the feasibility of large scale distributed digital libraries. The primary question we wish to ask is: “How many digital

library servers can be incorporated into a distributed digital library and continue to provide service?” To attempt to answer this question, a study to determine the Internet load generated by a distributed digital library under different server configurations and user activity levels should be performed.

Some of the primary scalability parameters to be considered when looking at the effect of Internet load include the total number of digital library servers, the total number of library objects, the size of the digital library objects, the number of queries being processed by the servers, and the number of objects being published. For this study we will limit our analysis to include the total number of digital library sites, the total number of queries represented by active user counts, the size of the digital library objects, and network throughput.

## **4.2 Justification**

This study is being done to verify the assumption that developing large scale distributed digital libraries is the logical direction to proceed. There are many conflicting approaches to digital library development and disagreement on the future basic architecture issues (Gladney et al. 1994; Arms et al. 1995; Graham 1995; Griffiths and Kertis 1995; Lagoze et al. 1996). Reports from early digital library research projects (Crawford 1995; Maly et al. 1995; Schnase et al. 1994) show us the breath and depth of the research problems to be resolved and the many directions the research is taking. By analyzing a very large scale distributed digital library model we may be able to provide some substance to discussions that are sometimes based on speculation and assumptions.

### 4.3 Analysis and comparison of simulation studies

In the two simulation performance studies mentioned in Section 3.3, there are a number of distinct differences between them and between the Generic Distributed Digital Library (GDDL) model we are presenting. Tables 4 through 8 show the major features of each model and study.

**Table 4.** Primary goal of the studies

<b>Inquery</b>	<b>NCSTRL</b>	<b>GDDL</b>
Used to analyze performance issues of the prototype distributed information retrieval system based on Inquery.	Performance evaluation and tuning and conducting what-if analysis for different configurations of NCSTRL.	Study the effect of digital library scaling and GDDL architecture changes on network traffic.

**Table 5.** Model components defined

<b>Inquiry</b>	<b>NCSTRL</b>	<b>GDDL</b>
Connection Server	Top Level	TLUI, LUI - User Interface
Inquiry Server	Region	LS - Local Site
Clients	Dienst Server	IS - Index Server
	Lite Server	MS - Metadata Server
	Merged Index Server	DS - Data Server
	Central Index Server	I - Index
	Backup Server	M - Metadata
	User Population	D - Data
	Query	PR - Retriever
		PP - Publisher

**Table 6.** Measurements used in the studies

<b>Inquiry</b>	<b>NCSTRL</b>	<b>GDDL</b>
Query evaluation time	Query inter-generation time	Network throughput
Document retrieval time	Server response time to queries	Average index size
Summary retrieval time	Transmission time of request from one server to another	Average metadata size
Connection server time		Average data size
Time to merge results		
Network time		

**Table 7.** Differences in study implementations

<b>Inquiry</b>	<b>NCSTRL</b>	<b>GDDL</b>
Yacsim process simulation	Visual Simulation Environment (VSE)	Analytic model analysis



**Table 8.** Parameters varied in the studies

<b>Inquiry</b>	<b>NCSTRL</b>	<b>GDDL</b>
Number of users	unknown	Number of PP, PR
Document collections		Number of IS, MS, DS
Terms per query		Number of LS
Query term frequency		Number of I, M, D
User think time		
Answers returned		
Workload		

The (Balci et al. 1998) paper does not provide results that can be studied or evaluated. It appears to be a usable tool for the NCSTRL implementers to utilize but without seeing the actual visual simulation there is little to be gained from the paper.

The Inquiry study (Cahoon and McKinley 1997) provides extensive background information and discussion concerning the design of the simulation and a thorough discussion of the results are clearly demonstrated in discussion and tables. These results can also be studied and used as guidance concerning issues that are relevant in designing distributed digital libraries.

#### **4.4 Discussion**

This study differs from both examples above in that a software simulation has not been conducted. It also differs in that we are examining a digital library, as a generic architecture not tied to specific implementation constructs. The first step in this study is to define the Generic Distributed Digital Library (GDDL) and then to analyze the network activities typical in a broad sense. A more extensive study would include the development of a simulation based on this generic design. This GDDL study is broad and

only provides a gross analytic solution to the question being asked. Although a generic distributed digital library model has been defined, much more work should be done to provide better analysis and results.

## SECTION FIVE

### PROBLEM ANALYSIS

The primary focus of this study is to determine the feasibility of large scale distributed digital libraries as defined in Section 2.3. To facilitate this study it is useful to dissect the anatomy of a distributed digital library into its component parts and use those components to define the architecture of a generic distributed digital library. To feel confident that the generic distributed digital library (GDDL) that is defined using these components is correct, we have taken these generic components and demonstrated that they can also be used to represent the architecture of three currently available production digital libraries. Table 9 outlines the component names and primary functions.

#### **5.1 Generic model design**

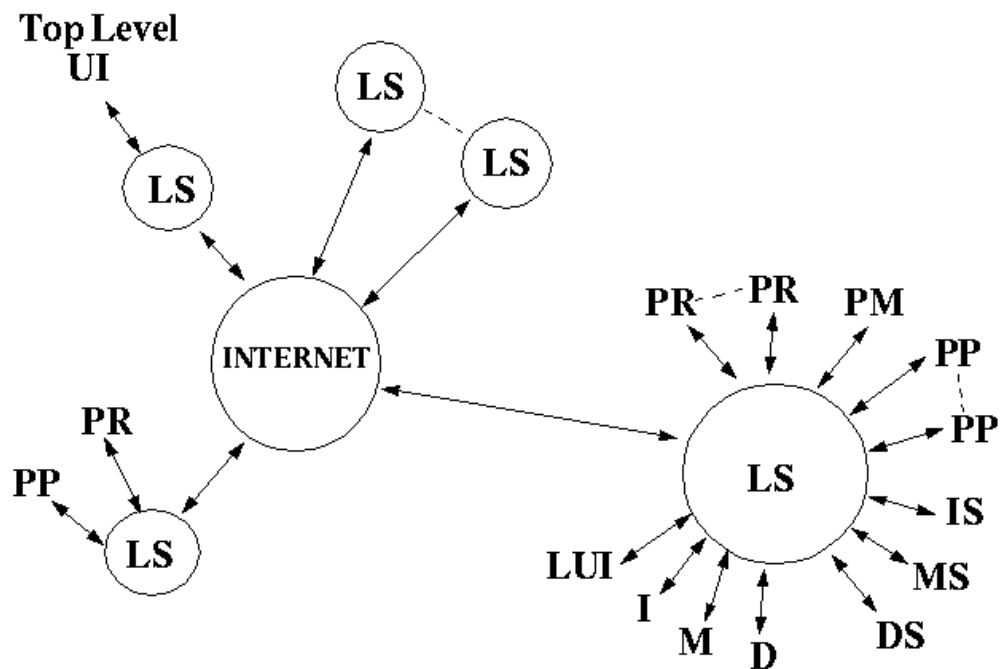
If we imagine a digital library as a set of independent objects serving unique functions with location independence then we could have a distributed digital library composed of data, metadata, and indices; the services that deliver this information; user interfaces and the people accessing these services.

##### **5.1.1 Description of system components**

Table 9 lists the basic components of a generic distributed digital library (GDDL). Shown in Figure 1 is a graphical representation of a GDDL. Definition 1 provides the basic definition of each of the components in the GDDL.

**Table 9.** Model nomenclature

Service Objects		People Objects	
TLUI	Top Level User Interface	PR	Retriever
LUI	Local User Interface	PP	Publisher
IS	Index Server	P M	Manager
MS	Metadata Server		
DS	Data Server		
I	Index		
M	Metadata		
D	Data		



**Figure 1.** Generic Distributed Digital Library model

**Definition 1.** The Generic Distributed Digital Library model components are:

*Internet* - The global networking infrastructure that interconnects the Local Sites.

*LS* - A Local Site can be single or multiple businesses, organizations, or entities connected via a local area network. In its simplest form a Local Site is a LAN for a single organization with one digital library in place for that organization.

*TLUI* - The Top Level User Interface provides search and retrieval access to all the Index, Metadata, and Data available at all the Local Sites. The Top Level User Interface can exist anywhere within the distributed digital library architecture.

*LUI* - The Local User Interface provides search and retrieval for the Local Site digital library Index, Metadata, and Data Servers.

*IS* - Index Server provides the service that accepts a request for index entries based on specified keywords for search. This service also creates, updates and manages the index. Each Local Site has at least one but possibly many Index Servers to manage indices of various collections of metadata and data.

*MS* - Metadata Servers provide access to synopsis information about the data as well as a high level view of the different representations of the data and supporting information.

*DS* - Data Servers provide mechanisms for the retriever to obtain the data in its various forms.

*I* - The Index object represents the body of indices being represented by the Index Servers.

*M* - The Metadata object represents the actual metadata information being maintained by the Metadata Servers.

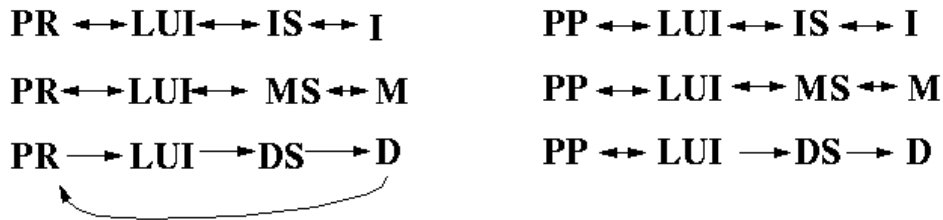
*D* - This is the Data Object.

*PP* - People publish into a Local Site digital library. The publish function is conducted by a user that has created a digital library object that includes Metadata and Data. These objects are inserted into the digital library through the Index Server, Metadata Server and Data Server.

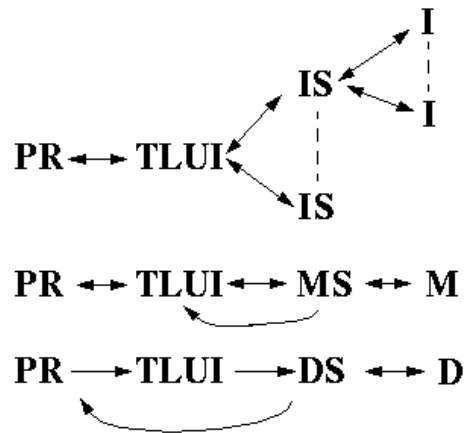
*PR* - People Retrieve represents the bulk of the day to day activities of the digital library. The People Retrieving can access the Top Level User Interface or any of the Local User Interfaces to search the Index, Metadata and Data at the Local Sites or across multiple sites throughout the distributed digital library.

### 5.1.2 Description of model data flow

In a distributed system, data of various kinds are constantly flowing in multiple directions. In a distributed digital library there are typical activities that occur with some regularity and in a somewhat predefined fashion. Represented in Figure 2 and Figure 3 are the data flow activities represented at a high level, expected to occur in the generic distributed digital library.



**Figure 2.** Local data flow of GDDL



**Figure 3.** Global data flow of GDDL

The basic activities have been separated into two categories, global and local. The Global activities are those things that go on at the Internet level. The local activities are occurring at the local site or Intranet level.

## 5.2 Model specifications

This model is designed to represent a generic distributed digital library and provide a basis for a future simulation implementation of this model. The model consist of multiple Local Sites distributed throughout the Internet and each Local Site may contain one or more Index, Metadata, and Data Servers; Index, Metadata and Data storage objects; a Local User Interface; and People Retrieving and People Publishing objects. There is one Internet and Top Level User Interface in this model and the location of the Top Level User Interface is arbitrary.

This model includes the components and parameters as shown in Table 10. By defining this model and parameters we not only see graphically the architecture of the GDDL, but also lay the framework for the development of a simulation for better analysis.



**Table 10.** GDDL model component specifications

<b>Object Identifier</b>	<b>Object Description</b>	<b>Object Parameters</b>
INET	Internet	INET ID TLUI ID TLUI Location (LS ID) Number of LS Number of Connections Size of Connections
TLUI	Top Level User Interface	TLUI ID LS ID
LS	Local Site	LS ID INET ID Number of IS, MS, DS, I, M, D, Number of UI, PR, PP, PM
LUI	Local User Interface	LUI ID LS ID TLUI ID
IS	Index Server	IS ID LS ID Number of I
MS	Metadata Server	MS ID LS ID Number of M
DS	Data Server	DS ID LS ID Number of D
I	Index	I ID IS ID LS ID Size
M	Metadata	M ID LS ID DS ID Size
PR	People Retrieve	LS ID Number of Queries Number of parameters Average Time
PP	People Publish	LS ID Number of M and D objects Size of objects

### 5.3 Examples of digital libraries

In various instantiations of digital libraries the independent service objects are often implemented in combination and tightly coupled by function and location. Although these objects exist in some form in the digital libraries being examined, their form takes many variations that have implications on performance, functionality, portability and maintainability. The examples shown in Table 11 represent this variety in Internet based digital library implementations.

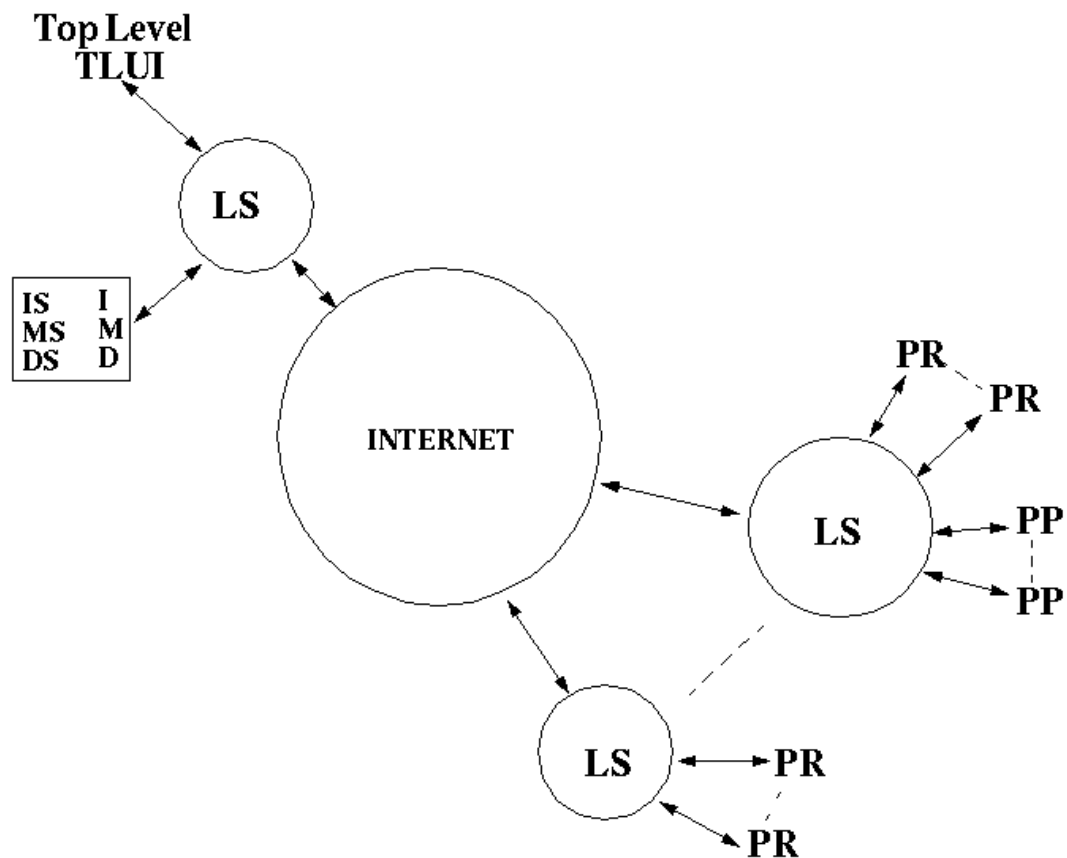
**Table 11.** Example digital libraries

<b>Digital Library</b>	<b>URL</b>	<b>Content</b>	<b># of Abstracts</b>	<b># of Reports</b>
Physics e-Print	<a href="http://xxx.lanl.gov">http://xxx.lanl.gov</a>	Physics and related technical papers	80 K	80 K
NTRS	<a href="http://techreports.larc.nasa.gov/cgi-bin/ntrs">http://techreports.larc.nasa.gov/cgi-bin/ntrs</a>	NASA technical reports	3.4 M	50 K
NCSTRL	<a href="http://www.ncstrl.org">http://www.ncstrl.org</a>	Computer Science technical reports	22 K	15 K

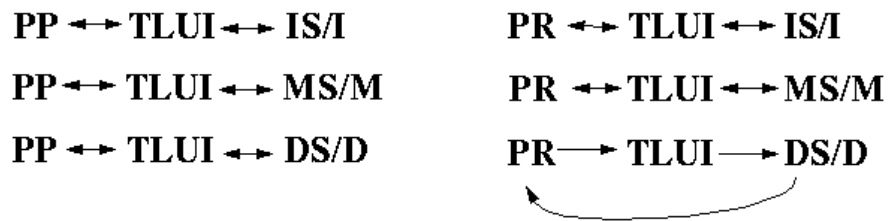
#### 5.3.1 Physics E-Print

The Physics E-Print digital library (Ginsbarg 1994) allows for remote Internet publisher and retriever access to the Index, Metadata and Data of the digital library through a Top Level User Interface that is tightly coupled with the Index, Metadata and Data Services. All the services provided by this digital library are implemented at a primary site with mirror sites providing duplicated service. Although this digital library

has a single primary site and is not truly a distributed digital library, it is included here because it provides distributed publishing, search and retrieval via the Internet. This also gives us a comparison model to visualize the difference in complexities between distributed and centralized digital library models.



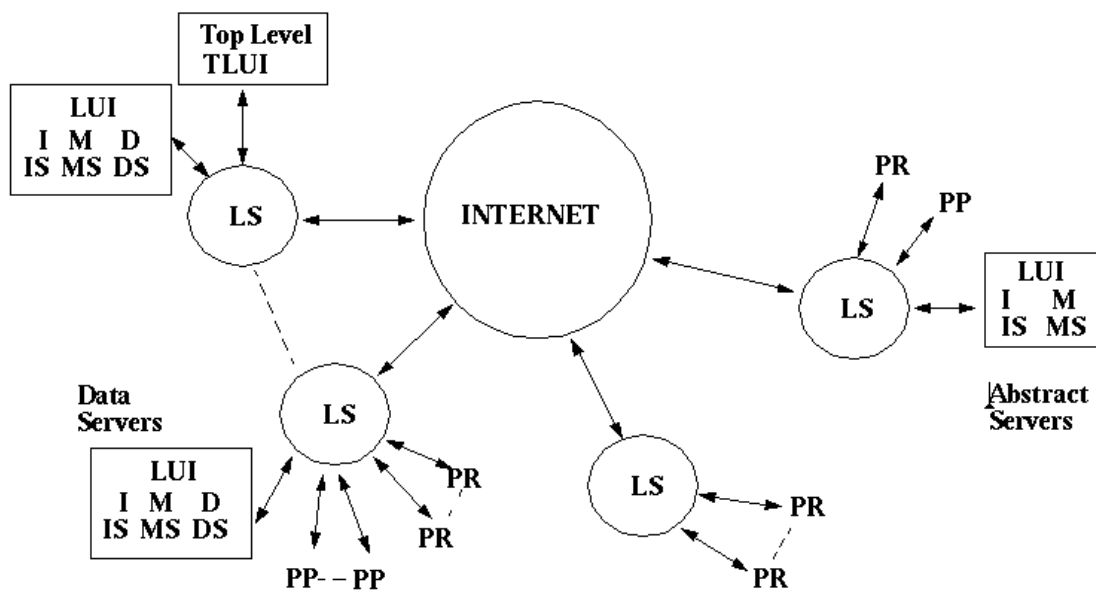
**Figure 4.** Physics E-Print model



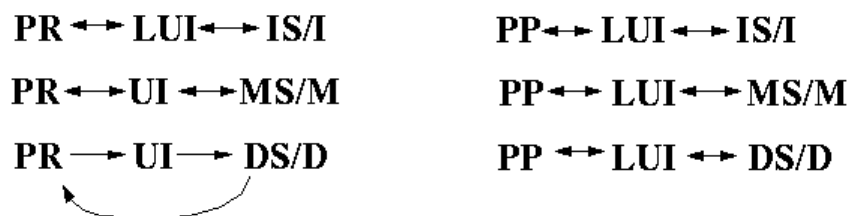
**Figure 5.** Global data flow of Physics E-Print

### 5.3.2 NTRS

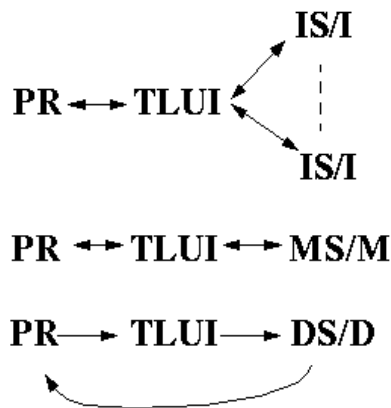
The NASA Technical Report Server (NTRS) digital library (Nelson et al. 1995) allows for local publishing and local and remote retrieving. The services are tightly coupled on single servers at each Local Site. There are 20 Local Sites distributed across the country and one Top Level User Interface site that provides search and retrieval access of all the Local Site information.



**Figure 6.** NTRS model



**Figure 7.** Local data flow of NTRS



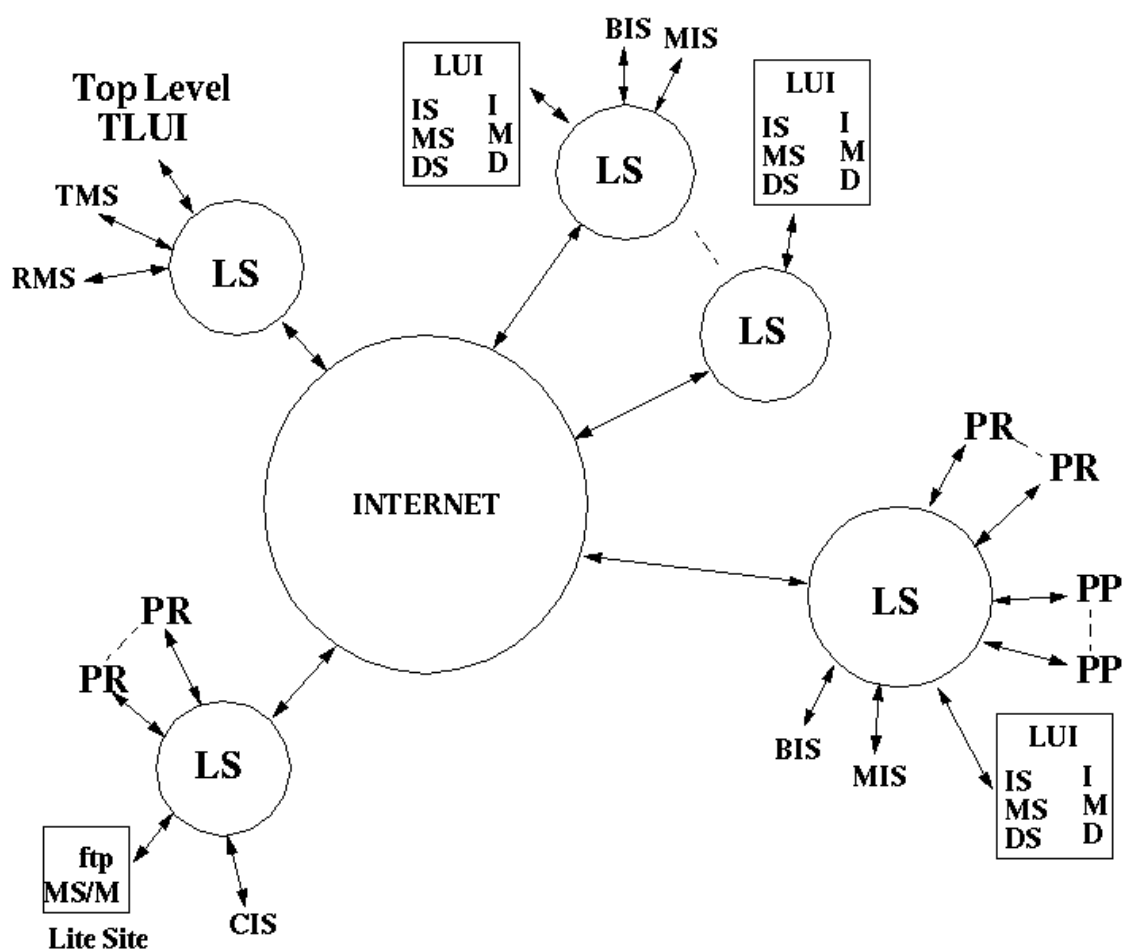
**Figure 8.** Global data flow of NTRS

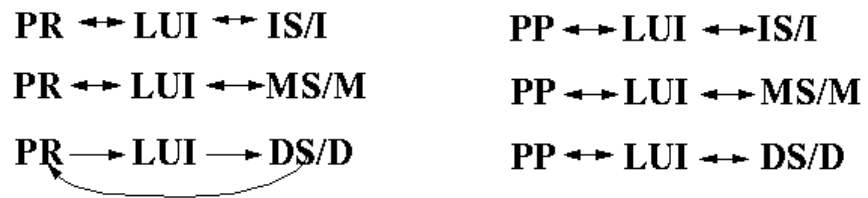
### 5.3.3 NCSTRL

The Networked Computer Science Technical Report Library (NCSTRL) is the most complex of the example digital libraries (Davis and Lagoze 1996) in this study. This digital library includes a Top Level User Interface, several Regional Sites and over 100 Local Sites. It also incorporates backup servers as well as top level and local Index and Metadata services. Because of the added complexity of this library we have defined additional components that are represented as a variation of the basic services provided in the generic distributed digital library. These additional components shown in Table 12 serve the same function as the Index Server and Metadata Server but at a higher level in the model.

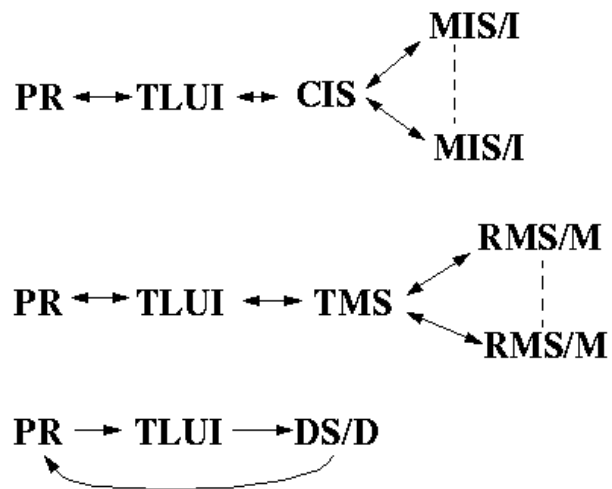
**Table 12.** NCSTRL specific components

<b>Object Name</b>	<b>Object Description</b>
CIS	Central Index Server
MIS	Merged Index Server
BIS	Backup Index Server
TMS	Top Metadata Server
RMS	Region Metadata Server

**Figure 9.** NCSTRL model



**Figure 10.** Local data flow of NCSTRL



**Figure 11.** Global data flow of NCSTRL



## 5.4 Parameters

Table 13 shows the list of parameters, as well as their description that will be considered in our analysis. For simplicity, we will assume at this time that a Local Site consists of one each of Index, Metadata, and Data Servers that serve one topic Index, the Metadata for this Index and the associated Data objects. It is expected that in a real world implementation, the number of the digital library components available at any given Local Site can vary greatly.

**Table 13.** Primary model parameters

<b>Parameters</b>	<b>Description</b>
# of PP	The number of People Publishing simultaneously
# of PR	The number of People Accessing the Digital Library for search or retrieval.
# of LS	The total number of Index, Metadata, and Data Servers in the DL.
# of IS, MS, DS	The total number of Local Sites in the Digital Library
# of I, M, D	The total number of Index, Metadata, and Data objects being served by the DL.
Size of I, M, D	The size in bytes of the Index, Metadata, and Data objects represented as an expected average.

## 5.5 Measurements and supporting data

It is important to understand current Internet technologies and future trends (Paxson 1997; Thompson et al. 1997) to evaluate the impact a distributed digital library architecture will have on the Internet. In Table 14 we show a variety of network technology and capacities available today as presented in (Tanenbaum 1996). Table 15

shows the throughputs that are being measured (Miller et al. 1998) for the vBNS high-performance network backbone (Jamison and Wilder 1997).

**Table 14.** Internet technology

<b>Technology</b>	<b>Gross Hardware Capacities</b>	<b>User Capacities</b>
OC 12	622.08 Mbps	445.824 Mbps
OC 3	155.52 Mbps	148.608 Mbps
OC1	51.84 Mbps	49.596 Mbps

**Table 15.** Internet throughputs

<b>Technology</b>	<b>Network Throughputs</b>	<b>Test Conducted</b>
OC 12	469 Mbps	UDP over ATM
OC 12	330 Mbps	TCP/IP over ATM
OC 3	130 Mbps	TCP/IP over ATM

The values shown in Table 16 were obtained from the Langley Technical Report Server (LTRS) (Nelson et al. 1994). LTRS is a subset of the NTRS and provides access to NASA Langley technical reports. The Indices in LTRS consist of a URL and title and do not vary considerably in size. Most metadata objects conform to a standard format and also have little variability in size. Data objects represent the greatest variability. The Data objects measured were all PDF files but they were generated from different original document formats including MS Word, PostScript and TIFF. The range of sizes represented in this average were from less than 40K to greater than 12MB.

**Table 16.** Average values measured from LTRS

<b>Data Object Name</b>	<b>Size in Bytes</b>
Index	468
Metadata	1,916
Data	1,457,389

It is important to note that the average size of 1.5 Mbytes is only representative for a digital library of text based technical reports. This number does not give any insight into the potential variability of size and types of data objects that can be made available and most likely will be a part of the digital library of the future. As such, it is probably a conservative number considering that the digital library of the future will be delivering video, audio, graphics, software, and large volume works such as books and data files.

We saw in Table 11 example digital libraries in use today. They range in size from 15,000 to 80,000 thousand reports, and anywhere from 22,000 to 3.4 million abstracts. We also know that these digital libraries are growing yearly. The Physic e-Print service reports they receive 18,000 new submissions yearly. The submission rates may grow, as the user communities better understand digital library technology and efficient means are provided to facilitate publishing into the libraries. As a digital library grows, so does the index to the volume of information. An important issue here is the time it takes to search an index is proportional to its size. In (Esler and Nelson 1997) we see a wide variety in the performance of index engines in part due to the size of the index being searched. This

has a direct affect on the response time users experience when searching a digital library. As the number of objects indexed in a digital library increases the overall performance of the digital library is expected to decrease after some critical point is reached. By distributing digital libraries as multiple smaller entities, this performance problem may be avoided.

## **5.6 Discussion of study cases**

Nine different cases represent the expected range of activities that occur at any point in time for a typical operational digital library. These activities occur concurrently and all contribute to the network traffic and load at the Local Site and on the Internet as well as to the load on the User Interface system and the different digital library service servers.

A typical user session will consist of a combination of searches and retrieval operations over a period of time with a great deal of intermixing of Index searches, Metadata retrievals and fewer Data retrievals. This general activity is represented in Table 17 and then broken down into smaller cases of activities shown in Table 18. The list of user actions includes the identifier for the Service object and People object active for each step in the session. The user connects to the TLUI, conducts a search of Index Services, retrieves Metadata, retrieves Data, and continues with these activities in an unpredictable way.

**Table 17.** User session characteristics

<b>User Action</b>	<b>Network Activity</b>
PR from TLUI	one to one
Search all LS/IS	one to many search
Return all I hits	many to one response
Retrieve one M from LS/MS/M	one to one
Retrieve one D from LS/DS/D	one to one
...reiterate between M and D...	
...reiterate from beginning...	

The nine cases shown in Table 18 are subdivided as either global or local based on the network traffic generated. Cases I through V represent global activities and generated Internet traffic while Cases VI through IX represent local activities and generate local traffic only.

**Table 18.** Digital library usage cases

<b>Characterization of User Action</b>	<b>Network Traffic</b>
<b>Case I – Global Query of all Index Servers</b>	<i>Internet traffic</i>
PR from TLUI	one to one
Search all LS/IS/I	one to many search
Return all I hits	many to one response
<b>Case II – Global Query of a Subset of Index Servers</b>	<i>Internet traffic</i>
PR from TLUI	one to one
Search some LS/IS/I	one to many search
Return all I hits	many to one response
<b>Case III - Global Query of one Index Server</b>	<i>Internet traffic</i>
PR from TLUI	one to one
Search one LS/IS/I	one to one search
Return all I hits	one to one response
<b>Case IV - Global Retrieval of Metadata</b>	<i>Internet traffic</i>
PR from TLUI	one to one
Request M	one to one
M transferred from LS/MS/M	one to one
<b>Case V - Global Retrieval of Data</b>	<i>Internet traffic</i>
PR from TLUI	one to one
Request D	one to one
D transferred from LS/DS/D	one to one
<b>Case VI - Local Site Publishing of Data</b>	<i>local site traffic</i>
PP to LUI	one to one
Submit one LS/IS/I	one to one submission
Confirmation returned	one to one response
<b>Case VII - Local Site Index Search</b>	<i>local site traffic</i>
PR from LUI	one to one
Search one LS/IS/I	one to one search
Return all I hits	one to one response
<b>Case VIII - Local Site Metadata Retrieval</b>	<i>local site traffic</i>
PR from LUI	one to one
Request M	one to one
M transferred from LS/MS/M	one to one
<b>Case IX - Local Site Data Retrieval</b>	<i>local site traffic</i>
PR from LUI	one to one
Request D	one to one
D transferred from LS/DS/D	one to one

### 5.6.1 Breakdown of cases studied

In this study we will assume from anecdotal evidence certain usage patterns for a typical digital library session. It is expected that users will at least spend part of the session in Case I, IV and V; index search, metadata retrieval and data retrieval. They may also spend time in Case II and III and a typical user session will have numerous metadata retrievals and fewer data retrievals. Given this, a session is suggested to have the percent values listed in Table 19.

Given this partitioning of a user session we can establish how many users to expect to be generating traffic based on case activity. For example, if we assume that we have 1000 simultaneous users, then the breakdown of activities will be as shown in Table 20. We can then use these numbers to calculate traffic generated per case for a given point in time and user population.

**Table 19.** Case breakdown by percentages

<b>Case</b>	<b>% of Time Sample A</b>	<b>% of time Sample B</b>
Case I	50	15
Case II	10	15
Case III	5	20
Case IV	20	35
Case V	15	15

**Table 20.** Case breakdown by user count

<b>Case</b>	<b>Number of Users Sample A</b>	<b>Number of Users Sample B</b>
Case I	500	150
Case II	100	150
Case III	50	200
Case IV	200	350
Case V	150	150

### 5.6.2 Case analysis

The traffic for each case has two directions. First the data going to the services in the form of requests being made for indices, metadata and data and then the data being returned to the user in the form of a list of indices, the metadata and the data objects. Some of this data flows from the user to the top-level user interface (TLUI) and then to the individual services and data also flows back to the TLUI for presentation to the user. Data objects are returned directly to the user and not routed through the TLUI.

In these formulas we are only interested in data being returned from the services to either the TLUI or directly to the user. We will not consider the traffic generated by the request for service from the User Interfaces. It is assumed that the amount of traffic generated by the user query and the User Interface search is less important compared to the total volume of data being returned to the user and the user interface. In this study we are only going to consider the traffic generated by the global cases and we are not distinguishing between traffic returning to the user or the user interface server. It is assumed that all the returning traffic must traverse the Internet and that is the number we are trying to establish.



Case VI, Case VII, Case VIII, and Case IX are not calculated because they represent local traffic only and do not have an impact on the total Internet traffic generated. Cases I through V are basic search and retrieval operations. The publishing activity is represented in Case VI and is considered a local activity based on the assumption that in most cases publishing is done at the users local site. We do expect some publishing to occur at the global level but we do not know at this time what percentage of all publishing will occur globally. We will assume this is a small enough percentage to not warrant inclusion in this study.

The formulas for Case I and II are a function of the total number of local sites being considered in the architecture multiplied by the worst case expected response of 250 indices returned per local site and the average number of bytes per indices.

Case III is the average indices size multiplied by the worst case number of responses. Case IV and V are assigned the values calculated from LTRS log data. No additional overhead is added to these numbers. The value for Case V has been rounded up for ease in calculation.

TLS represents the Total number of Local Sites to be varied in the study and T is used to represent the Total number of bytes generated per case. Table 21 shows the equations used to calculate the traffic generated per usage case for traffic returning to the user from the service.

**Table 21.** Equations used for case analysis

<b>Case</b>	<b>Worst Case</b>	<b>Average Case</b>
Case I	$T = (\text{TLS})(250)(468)$	$T = (\text{TLS})(10)(468)$
Case II	$T = .5(\text{TLS})(250)(468)$	$T = .5(\text{TLS})(10)(468)$
Case III	$T = (250)(468)$	$T = (10)(468)$
Case IV	$T = 1916 \text{ bytes}$	$T = 1916 \text{ bytes}$
Case V	$T = 1.5 \text{ Mbytes}$	$T = 1.5 \text{ Mbytes}$

### 5.7 Study assumptions

In a fully functional distributed digital library all activity, either local or global, has an impact on the total system performance. Because in this study we are focusing on Internet traffic generated, the traffic generated by local functions such as publishing and local queries will not be factored in. This assumes that people publish into the digital library at their Local Site and no Internet traffic is generated. It is reasonable to expect that in a real world distributed digital library, publishing may occur from any point in the system but it is also assumed that the level of this activity is insignificant and will be of little use in this analysis.

For the generic distributed digital library we are assuming all Local Sites are equivalent and all index are considered equal. Metadata and Data sizes are also considered equal and the averages presented are based on a Scientific Technical Information (STI) model. The byte counts were obtained from the NASA Langley Technical Report Server (LTRS) (Nelson et al. 1994) implementation through measurements and averaging of

existing contents. In a more diverse digital library each Local Site would vary greatly from the other Local Sites in total quantity, size and type of Index, Metadata, and Data being served.

Some Assumptions are presented for the digital library usage cases shown in Table 18. Case II assumes that a query of a subset of all Index Servers available would search 50% of these servers. This number could actually vary from the minimum of one represented in Case III to any number in between to the maximum number of Local Sites available represented in Case I. Assuming 50% is an attempt to capture the average. It is unknown what is typical in the situation when users are preselecting search sites. They may be selecting sites based on geographical, political, subject or personal preferences. This is an unknown factor to this author. Methods for reducing the number of servers queried is a subject of research (French et al. 1998). It is important to limit the number of servers searched too only those that can satisfy the query. This reduces the total network traffic and query processing time and results in a more efficient system.

In Cases I, II, and III of Index searches, we assume that 250 indices hits will be returned per Local Site Index Server queried. This represents the maximum allowable hits for a typical search engine configuration and is considered a worst case example. The logic behind this assumption is that there are no measured data available to show the typical number of indices returned per a global search. Even with data to examine concerning search hits and misses the characteristics vary so much that an average would not be a useful measure.

The subject of user query characteristics is broad and requires gathering large amounts of data related to user query analysis and system usability factors. In this study broad assumptions have been made concerning user query characteristics based on personal experience and anecdotal evidence. Further research and data gathering and analysis is needed to better define this aspect of the study.

In the calculation for total generated traffic we have to make some assumptions concerning how many active users there will be and what are the activities they are performing. We have defined nine different Cases of typical Digital Library activities but there is no data to tell us how many users are simultaneously interfacing with the digital library and what activities they are performing at any given time. Without doing a great deal of research into user usage patterns and system usage statistics we will assume a typical user usage pattern based on personal experience and make assumptions on the total user population counts.

## **5.8 Analytic formulas**

In a broad look at Internet Traffic we can say that the total Internet load created by a distributed digital library is minimally a function of the items shown in Equation 1. The query activities can be further broken down into more distinct parts as shown in Equation 2. To calculate the total Internet traffic generated from the services using the cases defined in Table 19 we use Equation 3.

**Equation 1.** Total Internet load

Total Internet Load =

$$\begin{aligned}
 &\text{All Global Query Activities (Case I, II, III)} && + \\
 &\text{All Global Publisher Activities (none)} && + \\
 &\text{All Global Metadata Retrieval Activities (Case IV)} && + \\
 &\text{All Global Data Retrieval Activities (Case V)}
 \end{aligned}$$

**Equation 2.** Global query Internet load

Global Query Internet Load =

$$\begin{aligned}
 &\text{Queries of all Index Servers (Case I)} && + \\
 &\text{Queries of a subset of Index Servers (Case II)} && + \\
 &\text{Queries of one Index Server (Case III)}
 \end{aligned}$$

**Equation 3.** Total Internet traffic

Total Internet Traffic =

$$\begin{aligned}
 &(\# \text{ of Case I})(\text{Case I traffic}) && + \\
 &(\# \text{ of Case II})(\text{Case II traffic}) && + \\
 &(\# \text{ of Case III})(\text{Case III traffic}) && + \\
 &(\# \text{ of Case IV})(\text{Case IV traffic}) && + \\
 &(\# \text{ of Case V})(\text{Case V traffic})
 \end{aligned}$$

## 5.9 Results tables

Traffic is defined as the total number of bytes that cross the Internet from the service through the TLUI or to the user for each case presented. Tables 22 and 23 show the calculation of traffic generated for each case as the total number of Local Sites is increased. Tables 24 through 27 show the calculation of final Internet traffic generated for different combinations of total number of users, cases, number of Local Sites, and sample usage patterns. The numbers for total users represents approximately 100 users accessing the Top Level User Interface per Local Site. This is a worst case analysis and the choice of 100 users is an arbitrary best guess based on the assumption that a Local Site represents some large organization or entity and that 100 users accessing the digital library at peak is reasonable to expect. This assumption is consistent with the expected growth patterns for the University of Illinois Digital Library Initiative as stated in (McGrath 1996).

Four different calculations were done to examine the worst case and average case results using two sample sets of user usage patterns show in Table 19. The worst case is determined by the use of 250 return indices per Local Site searched. The average case is determined by reducing the number of indices returned per Local Site to 10. Tables 23 and 24 show the amount of traffic generated for each individual case as defined by the equations shown in Table 21.

**Table 22.** Total traffic generated per individual case for worst case analysis

<b>Case</b>	<b>10 LS</b>	<b>100 LS</b>	<b>1,000 LS</b>	<b>10,000 LS</b>
Case I	1.17 Mb	11.7 Mb	117 Mb	1.17 Gb
Case II	.585 Mb	5.85 Mb	58.5 Mb	585 Mb
Case III	.117 Mb	.117 Mb	.117 Mb	.117 Mb
Case IV	1916 bytes	1916 bytes	1916 bytes	1916 bytes
Case V	1.5 Mb	1.5 Mb	1.5 Mb	1.5 Mb

**Table 23.** Total traffic generated per individual case for average analysis

<b>Case</b>	<b>10 LS</b>	<b>100 LS</b>	<b>1,000 LS</b>	<b>10,000 LS</b>
Case I	46,800 bytes	.468 Mb	4.68 Mb	46.8 Mb
Case II	23,400 bytes	.234 Mb	2.34 Mb	23.4 Mb
Case III	4680 bytes	4680 bytes	4680 bytes	4680 bytes
Case IV	1916 bytes	1916 bytes	1916 bytes	1916 bytes
Case V	1.5 Mb	1.5 Mb	1.5 Mb	1.5 Mb

**Table 24.** Calculation of total traffic for Worst Case using Sample A

<b>Total LS</b>	<b>Case I</b>	<b>Case II</b>	<b>Case III</b>	<b>Case IV</b>	<b>Case V</b>	<b>Total Users</b>	<b>Total Traffic</b>
10	500	100	50	200	150	1,000	874 MB
100	5,000	1,000	500	2,000	1,500	10,000	66 Gb
1,000	50,000	10,000	5,000	20,000	15,000	100,000	6 Tb
10,000	500,000	100,000	50,000	200,000	150,000	1,000,000	643 Tb

**Table 25.** Calculation of total traffic for Average Case using Sample A

<b>Total LS</b>	<b>Case I</b>	<b>Case II</b>	<b>Case III</b>	<b>Case IV</b>	<b>Case V</b>	<b>Total Users</b>	<b>Total Traffic</b>
10	500	100	50	200	150	1,000	251 Mb
100	5,000	1,000	500	2,000	1,500	10,000	5 Gb
1,000	50,000	10,000	5,000	20,000	15,000	100,000	282 Gb
10,000	500,000	100,000	50,000	200,000	150,000	1,000,000	25 Tb



**Table 26.** Calculation of total traffic for Worst Case using Sample B

<b>Total LS</b>	<b>Case I</b>	<b>Case II</b>	<b>Case III</b>	<b>Case IV</b>	<b>Case V</b>	<b>Total Users</b>	<b>Total Traffic</b>
10	150	150	200	350	150	1,000	512 Mb
100	1,500	1,500	2,000	3,500	1,500	10,000	29 Gb
1,000	15,000	15,000	20,000	35,000	15,000	100,000	3 Tb
10,000	150,000	150,000	200,000	350,000	150,000	1 M	263 Tb

**Table 27.** Calculation of total traffic for Average Case using Sample B

<b>Total LS</b>	<b>Case I</b>	<b>Case II</b>	<b>Case III</b>	<b>Case IV</b>	<b>Case V</b>	<b>Total Users</b>	<b>Total Traffic</b>
10	150	150	200	350	150	1,000	237 Mb
100	1,500	1,500	2,000	3,500	1,500	10,000	3 Gb
1,000	15,000	15,000	20,000	35,000	15,000	100,000	128 Gb
10,000	150,000	150,000	200,000	350,000	150,000	1 M	10 Tb

## **SECTION SIX**

### **FINDINGS AND INTERPRETATIONS**

What would happen if we introduced a 10,000 Local Site digital library onto the existing Internet? The highest capacity backbones currently available on the Internet range from 150 Mbps to 622 Mbps. Given this and some additional information we can calculate the approximate amount of time it would take to transfer data for the architecture examples calculated in Tables 24 through 27. Shown in Tables 28 and 29 is the approximate amount of time it would take to transfer the calculated amount of data for the four different situations represented as number of Local Sites in the GDDL. The values shown in Table 28 assume a network throughput of 130 Mbps. The values shown in Table 29 assume a network throughput of 450 Mbps. These throughput numbers were obtained from the vBNS web site and (Miller et al. 1998) and represent the capabilities of a finely tuned high-performance network.

As stated in Section 5.6.2, the total traffic numbers shown in Tables 28 and 29 are not broken down by destination. These numbers represent the traffic that we expect to cross an Internet backbone to various destinations. Additional useful information would be the percentage of this traffic that is returning to the Top Level User Interface and the percentage being dispersed directly to users distributed throughout the Internet. This would be helpful in determining the worst case expected load on the Top Level User Interface the GDDL.

**Table 28.** Time to transmit at 130 Mbps

	10 LS	100 LS	1,000 LS	10,000 LS
Worst Case/Sample A				
Total Traffic	874 Mb	66 Gb	6 Tb	643 Tb
Seconds	7	513	49,615	4,951,778
Worst Case/Sample B				
Total Traffic	512 Mb	29 Gb	3 Tb	263 Tb
Seconds	4	222	20,462	2,026,915
Average Case/Sample A				
Total Traffic	251 Mb	5 Gb	282 Gb	26 Tb
Seconds	2	37	2,171	199,769
Average Case/Sample B				
Total Traffic	237 Mb	3 Gb	128 Gb	10 Tb
Seconds	2	26	984	82,769

**Table 29.** Time to transmit at 450 Mbps

	10 LS	100 LS	1,000 LS	10,000 LS
Worst Case/Sample A				
Total Traffic	874 Mb	66 Gb	6 Tb	643 Tb
Seconds	2	148	14,333	1,430,513
Worst Case/Sample B				
Total Traffic	512 Mb	29 Gb	2 Tb	263 Tb
Seconds	1	64	5,911	585,553
Average Case/Sample A				
Total Traffic	251 Mb	5 Gb	282 Gb	26 Tb
Seconds	1	11	627	57,711
Average Case/Sample B				
Total Traffic	237 Mb	3 Gb	128 Gb	10 Tb
Seconds	1	8	284	23,911

Is it unreasonable to expect that there may be ten thousand local digital library sites distributed throughout the Internet at some point in the future? Or perhaps only one thousand digital libraries and if not, what will those local digital library sites consist of? Will the digital library of the future be supporting a small organization with a few users or will it be a digital library supporting a city, large business or government organization? It

is likely that much variety in digital library implementations will come forth supporting all quantities and types of data. We see many examples of this already (Crawford 1998). We also see evidence in (McGrath 1996) that it does seem reasonable to expect hundreds and even possibly thousands of so called digital libraries on the Internet of the future. There are many performance and scaling issues to consider for this to become a reality but without enough bandwidth all other issues become secondary.

Given the data in Tables 28 and 29, it is expected that a high-performance network infrastructure can support growth of distributed digital libraries well above 100 heavily used Local Sites but may have serious performance problems as it grew into the thousands. Beyond that, the problems of necessary bandwidth and other scaling issues become even more complex.

## **SECTION SEVEN**

### **FUTURE RESEARCH**

Much future research is possible on this topic. With the simple formulas and cases presented here more calculations and estimates can be made by varying the usage characteristics, the local site counts and the user counts. This would provide us with a range of possibilities from the low to high-end estimates of generated traffic results. Additional traffic calculations can also be made with data obtained from the three production examples presented (Physics e-Print, NTRS, and NCSTRL) and compared to the results obtained for the GDDL.

Additionally, the confidence in the results can be improved by eliminating many of the assumptions currently based on observation and anecdotal evidence. Specifically, it would be useful to obtain data concerning user usage patterns. This data could be obtained from current production digital library implementations if they can be set up to log the necessary data for analysis. The traffic data sizes used were narrowly defined by data obtained from one digital library implementation. A broad look at traffic patterns and sizes from a number of different types of digital libraries would provide a better average and more realistic results.

In this study we made assumptions concerning the definition of a local site. Because the field of digital libraries is young and examples are varied, we cannot say with any confidence what a local site will consist of. A further analysis of current digital

libraries and prototypes as well as World Wide Web patterns may yield some more insight into defining a digital library local site.

This study did not address the effect publishing or management functions may have on Internet traffic load. A better understanding of traffic routing patterns would also be useful to consider. It may reveal that not all the Internet traffic we are calculating is actually crossing the Internet. We do not fully understand the usage patterns to factor out queries and retrievals that logically cross the Internet but in actuality are local to the user.

We believe that there is enough complexity in this model and the problem analysis that a simulation of the model could be useful in providing better answers to the question asked. A simulation could also be beneficial in providing answers to as yet unasked questions concerning performance of the Top Level User Interface and Local Site activities.

Finally, there is much knowledge to be gained from the study of scaling and performance issues. The key will be in choosing the specific issues to study that will provide the most insight.

## SECTION EIGHT

### SUMMARY AND CONCLUSIONS

#### 8.1 Summary of contributions

In this study an attempt was made to define the low-level basic components of a generic distributed digital library and show that existing digital libraries do at least contain these components in some fashion. The purpose of this effort was to establish a basis for creating a generic digital library model for performance and scaling analysis separate and independent of any specific implementation issues found in performance studies done to date.

In addition to a Generic Distributed Digital Library (GDDL) definition, a set of user session Cases were defined that represent the primary distinct activities that users conduct when interfacing with a digital library. These cases were further differentiated based on type of network traffic generated, Intranet versus Internet. The cases that generated Internet traffic were further analyzed based on expected user activity level per case and this information was used to calculate expected internet traffic generated for a variety of user population counts and local site counts.

Finally, the information obtained from the case analysis and calculations of Internet traffic generated was used to determine the lower bound worst case analysis of future GDDL bandwidth needs. We see in Table 28 the time to transmit the calculated amount of data increases rapidly beyond 100 Local Sites at 130 Mbps throughput. In



Table 29 the time to transmit also increases but indicates that reasonable response could be expected with 100 to 500 Local Sites available.

## **8.2 Conclusions**

Due to the sheer volume of data potentially to be made available and the diversity in content and format it seems reasonable to suggest that a Digital Library network be established to provide information access service for all users to sites that conform to some digital library standards and capabilities. This would differentiate "Digital Libraries" from commercial, private and personal information sources and provide users reliable service to valid and sanctioned information for research, education and personal use. The expected communications needs are great and provide justification for this suggestion.

Access to a universal digital library that will provide access to individual digital libraries should be a user service provided by the information superhighway. As stated in the report on technical challenges (Willemsen 1995) for the information superhighway, "...the superhighway should provide a "seamless" web of features and services to users, with thousands of systems and components interacting or operating in a way that is transparent to the user." A universal digital library could be one of the services provided, conforming to the standards established for distributed digital libraries.

## REFERENCES

- Adler, S., Berger, U., Bruggermann-Klien, A., Haber, C., Lamersdorf, W., Munke, M., Rucker, S., Spahn, and H.: Grey Literature and multiple collections in NCSTRL. University of Hamburg, Department of Computer Science, Doc-001, January 1998
- Andresen, D., Yang, T., Egecioglu, O., Ibarra, O., Smith, and T.: Scalability Issues for High Performance Digital Libraries on the Wold Wide Web. In: *Proceedings of the Third Forum on the Research and Technology Advances in Digital Libraries*, May 1996, pp 139-148
- Arms, W.: Key Architectural Issues in the Digital Library. Corporation for National Research Initiatives, February 1995, Available at <http://www.cnri.reston.va.us/home/cstr/arch/slides.html>
- Baclawski, K., Smith, J. E.: High-Performance, Distributed Information Retrieval. Northeastern University, College of Computer Science, January 1995, Available at [http:// www.ccs.neu.edu/home/kenb/key/highperf/hp.html](http://www.ccs.neu.edu/home/kenb/key/highperf/hp.html)
- Balci, O., Bertelrud, A., Esterbrook, C., Nance, R.: Visual Simulation Environment. In: *Proceedings of the 1998 Winter Simulation Conference*, IEEE, Piscataway, NJ, December 1998b
- Balci, O., Ulusarac, C., Shah, P., Fox, E.: A Library of Resuable Model Components for the Visual Simulation of the NCSTRL System. In: *Proceedings of the 1998 Winter Simulation Conference*, IEEE, Piscataway, NJ, To appear December 1998b
- Borbinha, J. L., Ferreira, J., Jorge, J., Delgado, J.: Networked Digital Libraries: the Concept and a Case Study. Position paper presented at the ACM SIGIR-97 Workshop on Networked Information Retrieval, Philadelphia, July 1997, Available at <http://ciir.cs.umass.edu/nir97/borbinha/html/jlbnir.html>
- Browne, S., Dongarra, J., Fox, G. C., Hawick, K., Kennedy, K., Stevens, R., Olson, R., Rowan, T.: Management of the NHSE - A Virtual Distributed Digital Library. In: *Proceedings of the Second International Conference on the Theory and Practice of Digital Libraries*, June 11-13, 1995, Austin, TX, pp 57-63
- Cahoon, B., McKinley, K. S.: Performance Analysis of Distributed Information Retrieval Architectures. UM-CS-1995-054, Department of Computer Science, University of Massachusetts, Amherst, MA, June 1995

- Cahoon, B., McKinley, K. S.: Performance Evaluation of a Distributed Architecture for Information Retrieval. In: *Proceedings of Nineteenth Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, Zurich, Switzerland, August 1996, pp 110-118
- Cahoon, B., McKinley, K. S.: Evaluating the Performance of Distributed Architectures for Information Retrieval using a Variety of Workloads. Department of Computer Science, University of Massachusetts, Amherst, MA, 1997
- Callan, J., P., Croft, W. B., Harding, S. M.: The INQUERY retrieval system. In: *Proceedings of the 3<sup>rd</sup> International Conference on Database and Expert System Applications*, Valencia, Spain, 1992
- Chen, S-S., Chien, Y-T., Griffin, S., Leiner, B., Neches, R., Lal, N.: Agency perspective on the Digital Library Initiative. NSF, ARPA, NASA, DLI, May 1996, Available at <http://computer.org/computer/dli/r50022/agencies.htm>
- Crawford, D., (ed): *Digital Libraries, Communications of the ACM*. Volume 38, Number 4, April 1995
- Crawford, D., (ed): *Digital Libraries: Global Scope, Unlimited Access, Communications of the ACM*. Volume 41, Number 4, April 1998
- Eichhorn, G.: The Digital Library of the Astrophysics Data System. *Astrophysics and Space Science* 247, nos. 1-2, 1997, pp 189-210
- Davis, J., Lagoze, C.: A Protocol and Server for a Distributed Digital Technical Report Library. Cornell University, April 1994
- Davis, J. R., Krafft, D. B., Lagoze, C.: Dienst: Building a Production Technical Report Server. In: *Advances in Digital Libraries*, Springer-Verlag, 1995, pp 211-222
- Davis, J., Lagoze, C.: The Network Computer Science Technical Report Library. Technical Report TR96-1595, Cornell University Computer Science, July 1996
- Esler, S., Nelson, M.: NASA Indexing Benchmarks: Evaluating Text Search Engines. In: *Journal of Computer and Network Applications*, vol. 20, no. 4, 1997, pp 339-353
- Esler, S., Nelson, M. L.: The Evolution of Scientific and Technical Information Distribution. *Journal of the American Society of Information Science*, 49(1), 1998, pp 82-91

- French, J.: NCSTRL notes: Some Performance Issues. Department of Computer Science, University of Virginia, January 1996
- French, J., Powell, A., Creighton, III, W. R.: Efficient Searching in Distributed Digital Libraries. In: *Proceedings of The Third ACM Conference on Digital Libraries*, Pittsburgh, PA, June 1998, pp 283-284
- Fuhr, N.: A Decision-Theoretic Approach to Database Selection in Networked IR. University of Dortmund, Dortmund, Germany, January 1997
- Ginsparg, P.: First Steps Towards Electronic Research Communication. *Computer in Physics*, 8, 1994, pp 333-341
- Gladney, H., Ahmed, Z., Ashany, R., Belkin, N., Fox, E., Zemankova, M.: Digital Library: Gross Structure and Requirements (Report from a Workshop). IBM Research Report RJ 9840, May 1994
- Graham, P.: Requirements for the Digital Research Library. Rutgers University Libraries. July, 1995, Available at <http://aultnis.rutgers.edu/texts/DRC.html>
- Gravano, L., Chang, K., Garcia-Molina, H., Lagoze, C., Paepcke, A.: STARTS, Stanford Protocol Proposal for Internet Retrieval and Search. CS-TR-97-1580, Digital Library Project, Stanford University, January 1997
- Gravano, L., Chang, K., Garcia-Molina, H., Paepcke, A.: STARTS: Stanford Proposal for Internet Meta-Searching. In: *Proceedings of the 1997 ACM SIGMOD International Conference On Management of Data*, 1997
- Griffiths, J-M., Kertis, K.: Access to Large Digital Libraries of Scientific Information Across Networks. Graduate School of Library and Information Science, The University of Tennessee, 1995
- Jamison, J., Wilder, R.: vBNS: The Internet Fast Lane for Research and Education. *IEEE Communications Magazine*, January 1997
- Kahn, R., Wilensky, R.: A Framework for Distributed Digital Object Services. cnri.dlib/tn95-01, May, 1995. Available at <http://www.cnri.reston.va.us/home/cstr/arch/k-w.html>
- Lagoze, C., Lynch, C., Daniel, R.: The Warwick Framework: A Container Architecture for Aggregating Sets of Metadata. Cornell University Computer Science Technical Report TR-96-1593, June 1996

- Lagoze, C., Ely, D.: Implementation Issues in an Open Architectural Framework for Digital Object Services. Cornell University Computer Science Technical Report, TR95-1540, June 1995
- Lagoze, C., Shaw, E., Davis, J. R., Krafft, D. B.: Dienst: Implementation Reference Manual, Cornell Computer Science Technical Report TR95-1514, 1995
- Lesk, M.: *Practical Digital Libraries: books, bytes, and bucks*. Morgan Kaufmann Publishers, San Francisco, California, 1997
- Lynch, C., Garcia-Molina, H. (eds): Interoperability, Scaling, and the Digital Libraries Research Agenda: A Report on the May 18-19, 1995 IITA Digital Libraries Workshop. August 1995
- Maa, M.-H., Esler, S., Nelson, M. L.: Lyceum: A Multi-Protocol Digital Library Gateway. NASA TM-112871, July 1997
- Maly, K., French, J., Fox, E., Selman, A.: Wide Area Technical Report Service – Technical Reports Online. *Communications of the ACM*, 38(4), 45, 1995
- McGrath, R.: UIUC DLI Project Scale-up: A Technical Evaluation. National Center for Supercomputing Applications. University of Illinois, Urbana-Champaign, December 15, 1996, Available at <http://www.ncsa.uiuc.edu/People/mcgrath/DLI/Scaling>.
- Miller, G. J., Thompson, K., Wilder, R.: Performance Measurement on the vBNS. In *Proceedings of the Interop '98 Engineering Conference*, Las Vegas, NV, May 1998
- Nelson, M. L., Gottlich, G. L., Bianco, D. J.: World Wide Web Implementation of the Langley Technical Report Server. NASA TM-109162, September 1994.
- Nelson, M.L., Gottlich, G. L., Bianco, D. J., Paulson, S. S., Binkley, R.L., Kellogg, Y. D., Beaumont, C. J., Schmunk, R. B. Kurtz, M. J., Accomazzi, A.: The NASA Technical Report Server. *Internet Research: Electronic Networking Applications and Policy*, 5(2), 1995, pp 25-36
- Nelson, M., Maa, M-H.: Optimizing the NASA Technical Report Server. In: *Internet Research: Electronic Network Applications and Policy*, vol. 6, no. 1, 1996, pp 64-70

- Nelson, M., Maly, K., Shen, S. N. T., Zubair, M.: NCSTRL+: Adding Multi-Discipline and Multi-Genre Support to the Dienst Protocol Using Clusters and Buckets. In: *Proceedings of IEEE Forum on Research and Technology Advances in Digital Libraries*, April 1998, pp 128-136
- Nelson, M.: Old Dominion University CS745 Class notes. 1998, Available at [http://www.cs.odu.edu/~nelso\\_m/cs745](http://www.cs.odu.edu/~nelso_m/cs745).
- Nurnberg, P., Furuta, R., Leggett, J., Marshall, C., Shipman III, F.: Digital Libraries: Issues and Architectures. In: *Proceedings of Second Annual Conference on the Theory and Practice of Digital Librarie*, June 1995
- Paxson, V.: Measurements and Analysis of End-to-End Internet Dynamics. Ph.D. Thesis, Computer Science Department, University of California, Berkeley, April 1997
- Schatz, B., Chen, H.: Building Large Scale Digital Libraries. *IEEE Computer*, 29(5), 1996, pp 22-26
- Schatz, B., Mischo, W., Cole, T., Hardin, J., Bishop, A., Chen, H.: Federating Diverse Collections of Scientific Literature. *IEEE Computer*, 29(5), 1996, pp 28-36
- Schnase, J., Leggett, J., Furuta, R., and Metcalfe, T. (eds): *Proceedings of the First Annual Conference on the Theory and Practice of Digital Libraries*. College Station, Texas, June 1994
- Tanenbaum, A.: *Computer Networks*, Prentice Hall PTR, Upper Saddle River, New Jersey, 1996
- Thompson, K., Miller, G., J., Wilder, R.: Wide-Area Internet Traffic Patterns and Characteristics. *IEEE Network*, Vol. 11, No. 6, November/December 1997
- VanHeyningen, M.: The Unified Computer Science Technical Report Index: Lessons in Indexing Diverse Resources. In: *Proceedings of the 2<sup>nd</sup> International World Wide Web Conference*, October 19-21, 1994, pp 535-543
- von Ofenheim, W. H. C., Heimerl, N. L., Binkley, R., Curry, M., Slater, R., Nolan, G., Griswold, T., Kovach, R., Corbin, B., Hewitt, R.: NASA Image eXchange (NIX). NASA/TM-1998-206925, February 1998
- Willemssen, J. (ed): Information Superhighway: An Overview of Technology Challenges. Chapter Report GAO/AIMD-95-23, United States General Accounting Office, January 1995

REPORT DOCUMENTATION PAGE			Form Approved OMB No. 0704-0188	
Public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302, and to the Office of Management and Budget, Paperwork Reduction Project (0704-0188), Washington, DC 20503.				
1. AGENCY USE ONLY (Leave blank)		2. REPORT DATE December 1998		3. REPORT TYPE AND DATES COVERED Technical Memorandum
4. TITLE AND SUBTITLE Architectural Optimization of Digital Libraries			5. FUNDING NUMBERS	
6. AUTHOR(S) Aileen O. Biser				
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES)  NASA Langley Research Center Hampton, VA 23681-2199			8. PERFORMING ORGANIZATION REPORT NUMBER  L-17790	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)  National Aeronautics and Space Administration Washington, DC 20546-0001			10. SPONSORING/MONITORING AGENCY REPORT NUMBER  NASA/TM-1998-208956	
11. SUPPLEMENTARY NOTES Also published as a MS Thesis for the Old Dominion University Computer Science Department.				
12a. DISTRIBUTION/AVAILABILITY STATEMENT Unclassified-Unlimited Subject Category 66      Distribution: Nonstandard Availability: NASA CASI (301) 621-0390			12b. DISTRIBUTION CODE	
13. ABSTRACT (Maximum 200 words) This work investigates performance and scaling issues relevant to large scale distributed digital libraries. Presently, performance and scaling studies focus on specific implementations of production or prototype digital libraries. Although useful information is gained to aid these designers and other researchers with insights to performance and scaling issues, the broader issues relevant to very large scale distributed libraries are not addressed. Specifically, no current studies look at the extreme or worst case possibilities in digital library implementations. A survey of digital library research issues is presented. Scaling and performance issues are mentioned frequently in the digital library literature but are generally not the focus of much of the current research. In this thesis a model for a Generic Distributed Digital Library (GDDL) and nine cases of typical user activities are defined. This model is used to facilitate some basic analysis of scaling issues. Specifically, the calculation of Internet traffic generated for different configurations of the study parameters and an estimate of the future bandwidth needed for a large scale distributed digital library implementation. This analysis demonstrates the potential impact a future distributed digital library implementation would have on the Internet traffic load and raises questions concerning the architecture decisions being made for future distributed digital library designs.				
14. SUBJECT TERMS Digital Library, Architecture, Performance, Simulation, Internet, Scaling, Distributed Model			15. NUMBER OF PAGES 79	
			16. PRICE CODE A05	
17. SECURITY CLASSIFICATION OF REPORT Unclassified	18. SECURITY CLASSIFICATION OF THIS PAGE Unclassified	19. SECURITY CLASSIFICATION OF ABSTRACT Unclassified	20. LIMITATION OF ABSTRACT	